

Text-Prompted Diffusion Null-Space Model for Zero-Shot Image Restoration

Haijing Zhang (haijing@stanford.edu)

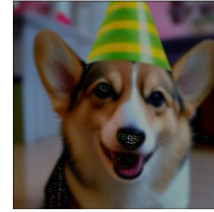
Yining Mao (yiningm@stanford.edu)



Colorization



Inpainting



Deblurring



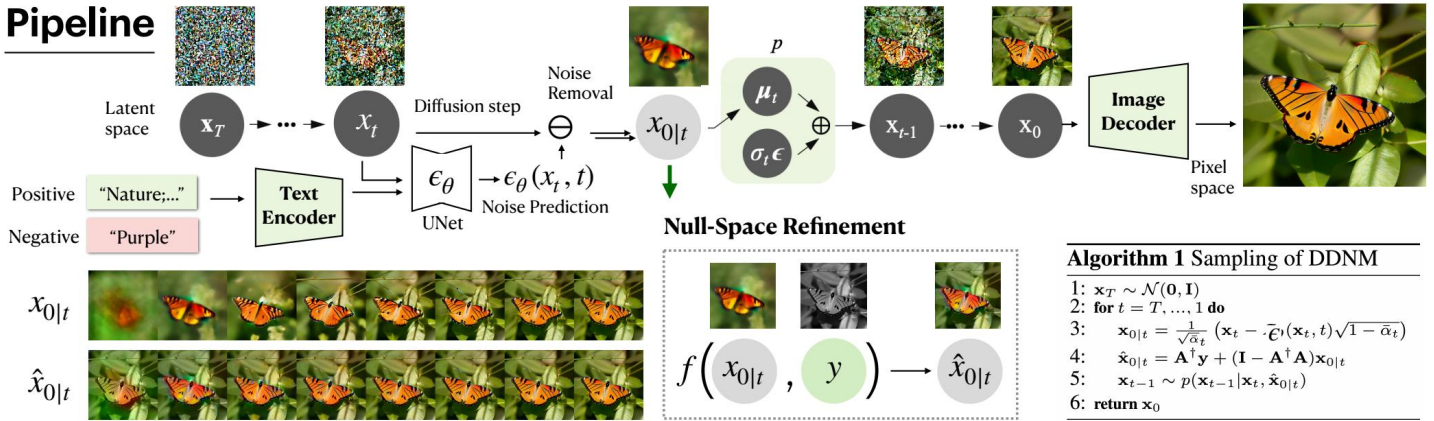
Abstract

[1] Yinhui Wang, Jiwen Yu, and Jian Zhang. Zero-shot image restoration using denoising diffusion null-space model, 2022.
 [2] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Essler, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition pages 10684–10695, 2022.

This project extends the usage of the **Denoising Diffusion Null-Space Model (DDNM)** [1], the state-of-the-art zero-shot image restoration method. With a pre-trained diffusion model, DDNM is capable of producing **realistic and data-consistent results without any training** or modifications to the network structure. However, undesirable results are generated occasionally due to the inherent randomness in the diffusion models and the performance limitations related to the pretrained model and training dataset.

To combat this problem, we propose integrating DDNM with the **text-prompt-enabled latent diffusion model** [2], widely known as Stable Diffusion. We use text prompts to fix the undesirable components in the result image, making the restoration process more **flexible and robust**. By reimplementing both algorithms and conducting experiments on three tasks, we have demonstrated that DDNM, originally developed for pixel space, is compatible with latent space diffusion models, and that using text prompts can effectively and efficiently make adjustments to the restoration results.

Pipeline

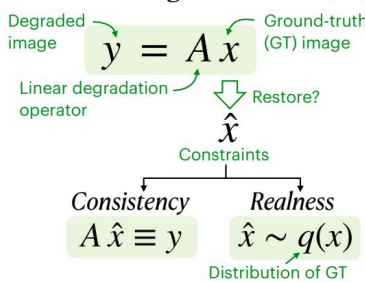


Algorithm 1 Sampling of DDNM

- 1: $x_T \sim \mathcal{N}(0, \mathbf{I})$
- 2: **for** $t = T, \dots, 1$ **do**
- 3: $x_{0|t} = \frac{1}{\sqrt{\alpha_t}} (x_t - \bar{c}_t(x_t, t)\sqrt{1 - \alpha_t})$
- 4: $\hat{x}_{0|t} = \mathbf{A}^\dagger y + (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A})x_{0|t}$
- 5: $x_{t-1} \sim p(x_{t-1}|x_t, \hat{x}_{0|t})$
- 6: **return** x_0

Denoising Diffusion Null-Space Model (DDNM)

Noise-free Image Restoration (IR)



Range-Null Space Decomposition

$$x \equiv \mathbf{A}^\dagger \mathbf{A}x + (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A})x$$

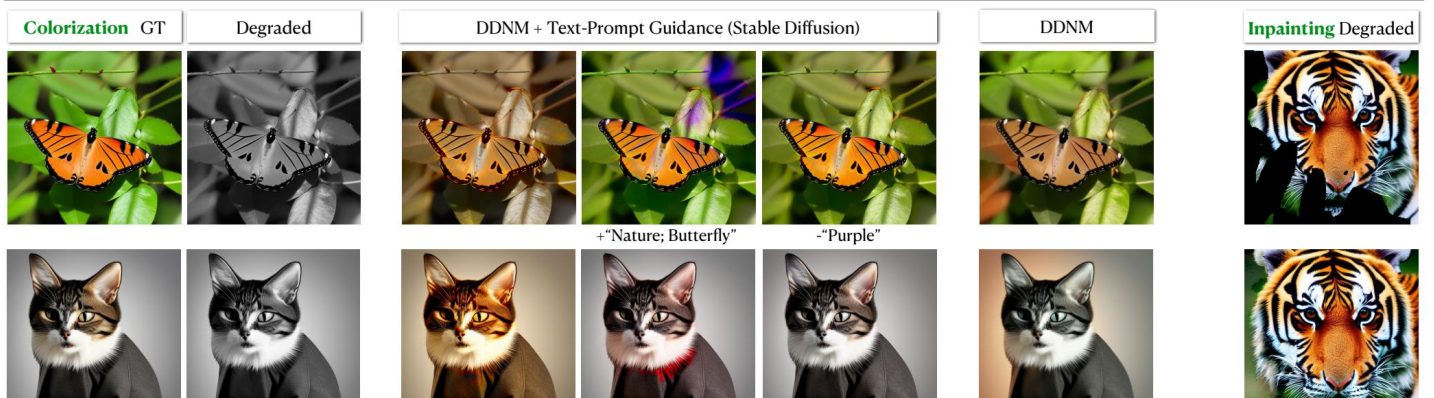
General solution $\hat{x} \equiv \mathbf{A}^\dagger y + (\mathbf{I} - \mathbf{A}^\dagger \mathbf{A})\bar{x}$

Multiply A on both sides $\mathbf{A}\hat{x} \equiv y$

This formula shows that no matter what \bar{x} is, it would always satisfy the data consistency. The quality of \bar{x} controls the realism.

For colorization:
 \mathbf{A} color \rightarrow grey
 \mathbf{A}^\dagger grey \rightarrow color

Experiments



Related Work

Traditional Model-based Method

Pros: less artifacts; Cons: fail in realistic details

$$\hat{x} = \arg \min_x \frac{1}{2\sigma^2} \|Ax - y\|_2^2 + \lambda R(x)$$

End-to-End Deep Learning Method

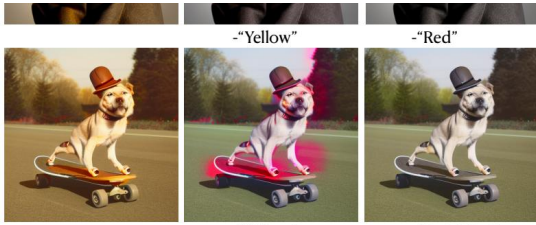
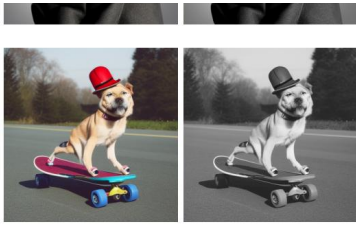
Pros: high-quality; Cons: task specific; lack interpretability

$$\hat{x} = \arg \min_{\theta} \sum_{i=1}^N \|D_{\theta}(y_i) - x_i\|_2^2$$

Zero-shot Pretrained Generative Model

Pros: good interpretability, efficient; Cons: hard to balance data consistency and realism

$$\hat{x} = \arg \min_w \frac{1}{2\sigma^2} \|AG(w) - y\|_2^2 + \lambda R(w)$$

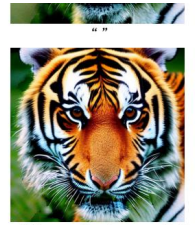
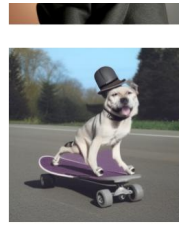


-“Yellow”

-“Red”

-“Yellow”

-“Red; Blue”



+“Tiger...; background...”